

Security and Cryptography 2019

Mirosław Kutylowski

IX. PRIVACY

Protection of personal data and GDPR

- declared as fundamental right in EU, but technically fundamental for cybersecurity
 - identity theft e.g. for financial criminality
 - mobbing, discrimination, social abuse
 - lack of protection is a threat for economy and national security

Frameworks

- EU and European Economic Area, adopted by many countries, some recognized as equivalent by EU
- Privacy Shield: <https://www.privacyshield.gov/>
- California Consumer Privacy Act

Privacy in communication

- one can hide payload communication
- it is not trivial how to hide **who is communicating with whom**
- this is a sensitive data

protection methods:

- broadcast channel
- token ring
- dining cryptographers -DC nets
- onion protocols and TOR

DC-nets

- protecting the source of a 1-bit message. One of cryptographers is sending a 0 or 1.
- protocol for 3 cryptographers sitting at a round table:
 - 1 each pair of neighbors establish a shared bit at random
 - 2 each cryptographer that is not transmitting computes XOR of the bits shared with the neighbors,
 - 3 the sender computes the same XOR but swaps it if the bit transmitted is 1
 - 4 each cryptographer reveals his result
 - 5 the message is the XOR of the bits published:
 - if the message is 0, then each shared bit occurs twice:

$$(b_{AB} \oplus b_{BC}) \oplus (b_{BC} \oplus b_{CA}) \oplus (b_{CA} \oplus b_{AB}) = 0$$

- otherwise, one of the bits is swapped: e.g. we have

$$(b_{AB} \oplus b_{BC}) \oplus (b_{BC} \oplus b_{CA} \oplus 1) \oplus (b_{CA} \oplus b_{AB}) = 1$$

Communication steganography

Idea: hiding data in innocent data transmitted (e.g. images, sound, protocol execution data)

steganography versus watermarking:

watermarking is not annoying but hard to remove, steganography is invisible

steganography in images

1. original picture taken – name: stego image
2. marking algorithm (maybe with a secret key) applied to message and the stego image
3. outcome transmitted/published
4. retrieving the covered message

invisibility: without a key impossible to decide whether there is a message hidden

typical applications: copyright protection, DRM,

Concrete techniques for image/video/audio steganography:

- changing LSB bits of gray scale
- JPEG encoding: cosinus transform, high frequency components are manipulated anyway for compression
- other digital transforms
- echoing
- audio encoding: transformation and assigning coefficients to waves – manipulations of certain coefficients undetectable for listeners

watermarking network flow:

- 1 based versus 0/1 based (in the first case the change is 1 no change encodes 0)
- all random parameters transmitted in clear may contain watermarks
- simple timing: interpacket delays, departure times
- mean balancing:
 - $2d$ probes divided into two sets A and B .
 - the expected values of A and B are the same with no watermarking
 - changing the characteristics so that expected values differ in some direction (the direction is the watermarked value)
- sources of mean balancing watermarks:
 - interpacket delays
 - interval centroid: divide into $2d$ intervals, in each compute mean arrival time
 - interval counting: divide into $2d$ intervals, in each compute the number of packets

- size: packet size (harder if block encryption applied), object size (https, malicious Javascript generating watermarked size data)
- network rate:
 - one can influence it with dummy packets
 - burst traffic (e.g. makes TOR useless in hiding routes)
- response times in transmission, packet order etc

Defense against watermarking network flow:

(sometimes problematic or illegal due to intellectual property rights)

- i. compression
- ii. transforms and random distortions
- iii. stretching (StirMark)
- iv. printing and scanning, input to an analog device and digitalize again

effectiveness measured by relative entropy:

$$D(P_1||P_2) = \sum_{q \in \text{space}} \text{Pr}_1(q) \cdot \log \frac{\text{Pr}_1(q)}{\text{Pr}_2(q)}$$

relative entropy of plaintexts and ciphertext should be smaller than some small ϵ

Mixer

messages m_1, m_2, \dots, m_k go through a mixer A :

- message m_i sent encrypted with the public key of A

$$C_i := \text{Enc}(\text{PK}_A, m_i)$$

- A decrypts C_1, \dots, C_k and forwards them to their destinations

Conditions:

encryption might be secure but nevertheless one can link the ciphertexts with the decrypted texts:

- message size
- timing

So:

- all messages should have the uniform size
- A should collect them all, decrypt and send them in a random order

Onion Routing

an attempt to hide the sender and the destination of a message – hiding in a crowd of messages

onion creation:

an onion created for a route going through servers A, B, C, \dots, Z to the final destination Γ

$$O_1 = \text{Enc}_A(B, \text{Enc}_B(C, \text{Enc}_C(\dots(\text{Enc}_Z(\Gamma, M))\dots)))$$

(by encryption Enc_X we mean encryption with the public key of X)

onion processing:

– O_1 sent from the origin machine to A ,

– server A decrypts with its private key and gets B and

$$O_2 = \text{Enc}_B(C, \text{Enc}_C(\dots(\text{Enc}_Z(\Gamma, M))\dots))$$

– A sends O_2 to B

– server B decrypts O_2 with its private key and gets C and $O_3 = \text{Enc}_C(\dots(\text{Enc}_Z(\Gamma, M))\dots)$

– the process is continued in the same way ...

– ... until server Z finds Γ, M and forwards the message M to machine Γ

each processing steps is like peeling off one layer of an onion

Limitations

- **idea**: if two or more onions enter a server at the same time, get partially decrypted, then forwarded, then it is impossible to say which incoming onion corresponds to which outgoing onion - **node mixing**
 - **traffic analysis**: assigning probabilities to permutations ($\pi(i) = j$ means that the i th sender has a message to the j th receiver)
 - **it is not enough to say that $\pi(i) = j$ with ppb $\approx \frac{1}{n}$** :
 - let us assume that the adversary knows that π is a circular shift
 - assume that the adversary gets extra knowledge:
 - "if source i is talking to destination j , then $i + 1$ is talking to destination $j + 1$ "
- however, still $\Pr(\pi(i) = j) = \frac{1}{n}$

Question: necessary length of the onions?

analytical results for restricted case of n senders and n receivers, messages sent simultaneously:

- $O(\log^2 n)$ if the adversary has a full knowledge of the system (not likely to have a better estimation unless ... big progress in math), **assumption**: uniform distribution for choosing destinations
- $O(\log n)$ if the adversary can see only a constant fraction of nodes, **assumption**: sender i may have non-uniform distribution of destination points
- it is easy to see that $\Omega(\log n)$ is necessary

meaning of the results:

traffic analysis does not improve our prior knowledge in a significant way (e.g. if we know in advance that source i always sends to destination j , then onions cannot hide this fact)

the guarantees are given in terms of **total variation distance** of two probability distributions:

$$\|\pi, \mu\| = \frac{1}{2} \sum_{\omega \in \Omega} |\pi(\omega) - \mu(\omega)|, \quad \text{where } \Omega \text{ is the set of all events}$$

Problems with onions

- **replay attacks**: just send the same onion (or partially decrypted onion) for the 2nd time: the same subonions will appear along the forwarding path

defense: universal reencryption, example based on ElGamal encryption:

- ciphertext of m :

$(\beta^r, m \cdot g^r, \beta^s, g^s)$ for r, s chosen at random

- reencryption:

1 choose α, β at random

2 replace (y_1, y_2, y_3, y_4) by $(y_1 \cdot y_3^\alpha, y_2 \cdot y_4^\alpha, y_3^\beta, y_4^\beta)$

thereby we get $(\beta^{r+\alpha s}, m \cdot g^{r+\alpha s}, \beta^{s\beta}, g^{s\beta})$

if order of the group is a prime number, then this is equivalent with choosing $(\beta^{r'}, m \cdot g^{r'}, \beta^{s'}, g^{s'})$ for random r', s'

- **local view**: not all users have the same list of servers

then: **long routes do not improve anonymity**. Toy example:

- give user A a list of servers with 50% of servers used by nobody else
 - no matter how long is the routing path designed by A , it is likely that close to destination the path goes through a rogue server
 - a few destinations available from this rogue server (50% of cases the rogue server sends directly to the destination)
 - an onion going through the rogue server originates from the attacked source
 - **timing at nodes**: delays necessary
- defense: collecting enough onions and flashing them at once. (**slowdown!!!**)
- **sparse traffic** means no protection

TOR

- free BSD licence
- connection based protocol, new connection established periodically (“10 minutes or so”)
- routes limited to 3 TOR nodes
- onion based forwarding the symmetric keys
 - i each node on the path learns only the predecessor and the successor
 - ii the path established step by step:
 - after establishing a subpath X_0, X_1, \dots, X_k the subpath is used to send an encrypted message over the channel to X_k stating that the next node is X_{k+1} .
 - the sender and X_{k+1} negotiate a new connection key via DH key exchange
 - iii after making a connection the message is encrypted symmetrically with the keys:

$$\text{AES}_{\text{relay1}}(\text{AES}_{\text{relay2}}(\text{AES}_{\text{relay3}}(m)))$$

each relay node removes one layer

- iv the messages back to the sender: instead of decryption: encryption with keys shared with the sender. The sender has to decrypt the onion

Problems:

- exit node knows the plaintext
- traffic correlation
- application level attacks
- Heartbleed - change of public keys, some clients use old keys,

Other issues:

- many authorities fight against TOR as it helps to escape the control

PSEUDONYMIZATION

Symmetric methods:

- **hashing the identifier:** $\text{pseudonym} = \text{Hash}(\text{identifier})$

problem: it is impossible to compute the identifier from the pseudonym, however hashing all possible identifiers and brute force reveals the link between the pseudonym and identifier

- **encryption with a (secret) symmetric key:** unlinkability, however the user cannot compute the pseudonym himself and the owner of the secret key can link all pseudonyms
- **hashing with a key:** as above, the party holding the secret key has to perform brute force to link back the pseudonym to the identifier

Asymmetric pseudonymization methods:

- based on Diffie Hellman Problem:

- a **domain** (service provider, database, etc) holds a pair of keys $(d, D = g^d)$
- a user Alice holds a pair $(x, X = g^x)$
- the **pseudonym** of Alice corresponding to D is $g^{x \cdot d}$, which is computed as X^d by the domain manager, and D^x by the user
- nobody but the user and the domain manager can compute the pseudonym:

for a 3rd person deciding whether X^d corresponds to X in domain D requires solving DDH Problem

- a variant based on domain and a central Authority:
 - the key d is not known to the domain authority
 - $d = d_A + d_{\text{domain}}$, where d_A is known by the authority and d_{domain} is known by the domain manager
 - steps of generating the pseudonym by the authority:
 - 1 the Authority computes $X' := X^{d_A}$ and presents X' to the domain manager
 - 2 the domain manager computes pseudonym as $(X')^{d_{\text{domain}}}$
 - linking a pseudonym with the starting public key is a reverse process but **both the domain manager and the Authority must participate** in it

- a variant from German personal identity cards (Restricted Identification):
 - **pseudonym** of a user with public key $X = g^x$ is $\text{Hash}(D^x)$
 - **pseudonym presentation**: by the ID card over a secure channel,
 - no proof that the pseudonym is correct
 - but a smart card can create only one pseudonym per domain
 - **revocation**: by computing $\text{Hash}((X^{d_A})^{d_{\text{domain}}})$ jointly by the Authority and the domain manager and putting the result on the **blacklist**
 - blacklisting a black sheep based on the domain pseudonym requires brute force and recomputing all pseudonyms
- more flexibility, if pairing groups are available but be careful: DDH might be easy and so the above methods do not work

Advantages and disadvantages of Restricted Identification:

- different pseudonyms generated automatically is
 - user friendly
 - makes re-identification based solely on data related to the pseudonym much harder
- problems:
 - converting a pseudonym in domain D_1 to a pseudonym in domain D_2 might be hard or infeasible, and require cooperation with the user and/or an authority
(problem area: moving pseudonymized medical records)

DATABASES and PRIVACY for QUERIES

the main problem is answering queries: does a query result disclose personal data?

Approach 1: **anonymity set**

- a query accepted if the number of record used to answer the query is at least k (and each concerns a different person)
- the method is naive: the attack is to ask for two sets of records: one including Alice and one excluding Alice to know the value for Alice

Approach 2: differential privacy

classify the algorithms (queries)

algorithm A satisfies ϵ -differential privacy, if for any two databases D and D' that differ by elimination of one record:

- for any subset S of the image of A :

$$\Pr(A(D) \in S) \leq e^\epsilon \cdot \Pr(A(D') \in S)$$

where the probability is over the random choices of A

Then:

- $\epsilon = 0$ is the ideal for privacy, but the result does not depend on the database contents (noise)
- so it is necessary to find balance between privacy (ϵ as small as possible) and information in the response (ϵ as big as possible)

Problem with outliers

with some records that have very different values it is hard to keep promise of *differential privacy*

solution: disregard them (as private data leak anyway) and concentrate on the rest

PSEUDONYMOUS SIGNATURES

Application areas:

- while having the pseudonyms, how to authenticate digital data? Digital signatures would solve the problem
- implementing GDPR rights in practice: a data subject can authenticate the request (e.g. for data rectification) in a database with pseudonyms by sending a request with a signature corresponding to the pseudonym

BSI Pseudonymous Signature:

- keys:

- domain parameters D_M and a pair of global keys (PK_M, SK_M)
- public key PK_{ICC} for a group of eIDAS tokens, the private key SK_{ICC} known to the issuer of eIDAS tokens
- assigning the private keys for a user:

the issuer chooses $SK_{ICC,2}$ at random, then computes $SK_{ICC,1}$ such that

$$SK_{ICC} = SK_{ICC,1} + SK_M \cdot SK_{ICC,2}$$

- a sector (domain) holds private key SK_{sector} and public key PK_{sector} .
- a sector has revocation private key $SK_{revocation}$ and public key $PK_{revocation}$
- sector specific identifiers $I_{ICC,1}^{sector}$ and $I_{ICC,2}^{sector}$ for the user:

$$I_{ICC,1}^{sector} = (PK_{sector})^{SK_{ICC,1}}$$

$$I_{ICC,2}^{sector} = (PK_{sector})^{SK_{ICC,2}}$$

- **signing:** with keys $SK_{ICC,1}$, $SK_{ICC,2}$ and $I_{ICC,1}^{\text{sector}}$ and $I_{ICC,2}^{\text{sector}}$ for PK_{sector} and message m

i choose K_1, K_2 at random

ii compute

- $Q_1 = g^{K_1} \cdot (PK_M)^{K_2}$

- $A_1 = (PK_{\text{sector}})^{K_1}$

- $A_2 = (PK_{\text{sector}})^{K_2}$

iii $c = \text{Hash}(Q_1, I_{ICC,1}^{\text{sector}}, A_1, I_{ICC,2}^{\text{sector}}, A_2, PK_{\text{sector}}, m)$

(variant parameters omitted here)

iv compute

- $s_1 = K_1 - c \cdot SK_{ICC,1}$

- $s_2 = K_2 - c \cdot SK_{ICC,2}$

v output (c, s_1, s_2)

- **verification:**

compute

- $Q_1 = (\text{PK}_{\text{ICC}})^c \cdot g^{s_1} \cdot (\text{PK}_M)^{s_2}$
- $A_1 = (I_{\text{ICC},1}^{\text{sector}})^c \cdot (\text{PK}_{\text{sector}})^{s_1}$
- $A_2 = (I_{\text{ICC},2}^{\text{sector}})^c \cdot (\text{PK}_{\text{sector}})^{s_2}$
- recompute c and check against the c from the signature

- why it works?

$$\begin{aligned}(\text{PK}_{\text{ICC}})^c \cdot g^{s_1} \cdot (\text{PK}_M)^{s_2} &= (\text{PK}_{\text{ICC}})^c \cdot g^{K_1} \cdot (\text{PK}_M)^{K_2} \cdot g^{-c \cdot \text{SK}_{\text{ICC},1}} \cdot (\text{PK}_M)^{c \cdot \text{SK}_{\text{ICC},2}} \\ &= (\text{PK}_{\text{ICC}})^c \cdot g^{K_1} \cdot (\text{PK}_M)^{K_2} \cdot g^{-c \cdot \text{SK}_{\text{ICC},1}} \cdot (g)^{-c \cdot \text{SK}_M \cdot \text{SK}_{\text{ICC},2}} \\ &= (\text{PK}_{\text{ICC}})^c \cdot g^{K_1} \cdot (\text{PK}_M)^{K_2} \cdot g^{-c \cdot \text{SK}_{\text{ICC}}} = g^{K_1} \cdot (\text{PK}_M)^{K_2} = Q_1\end{aligned}$$

- there is a version without A_1, A_2 and the pseudonyms $I_{\text{ICC},1}^{\text{sector}}, I_{\text{ICC},2}^{\text{sector}}$

- **Problems:**

- the authorities know the private keys (there is a way to solve it when the user gets two pairs of keys on the device and takes their linear combination)
- breaking into just 2 devices reveals the system keys
- possible to create a trapdoor for enabling to link pseudonyms
 - apart from $SK_{ICC} = SK_{ICC,1} + SK_M \cdot SK_{ICC,2}$ there is another relationship for the user u

$$x_u = SK_{ICC,1} + s_u \cdot SK_{ICC,2}$$

- x_u and s_u are dedicated for user u - maybe not in the database but derived from a secret key, say Z
- domain trapdoor: $T_{\text{domain},u} = PK_{\text{domain}}^{x_u}$ and s_u (it can be derived from Z alone)
- then one can conclude that nym_1 and nym_2 correspond to user u , if:

$$T_{\text{domain},u} = \text{nym}_1 \cdot \text{nym}_2^{s_u}$$

ANONYMOUS CREDENTIALS

two commercial products (libraries): Idemix (IBM) and UProve (Microsoft)

some details concerning Idemix

components:

- **actors:** issuer, recipient, verifier, trusted party
- **attributes:** for each attribute there is: name, value and type. The types are `int`, `string`, `date`, `enum` (enumeration). The attributes concern the recipient.
- **credentials:** given by the issuer to the recipient
 - i known (A_k): the issuer knows the value of an attribute
 - ii committed (A_c): the issuer knows a commitment to the attribute but not the commitment itself
 - iii hidden (A_h): the attribute is completely hidden to the issuer

– **keys:**

- single master key for each user (m_1)
- single master key for the Issuer – for creation of CL signatures

– **pseudonyms:**

- a single domain pseudonym for a user per domain: generated as as

$$\text{dom}^{m_1}$$

where dom is the public key of a domain, and m_1 is the user's master key

- pseudonyms are unlinkable

Cryptographic schemes used by Idemix

CL signatures:

- RSA group, special choice of primes: $p = 2p' + 1$, $q = 2q' + 1$, where p' and q' are primes
- choose at random quadratic residues: R_1, \dots, R_l, Z, S
- public key: $(n, R_1, \dots, R_l, Z, S)$, private key: p, q (enabling computation of roots mod n)
- security based on **Strong RSA assumption**: it is infeasible to compute e -roots for $e > 2$
- signature for messages m_1, \dots, m_l :
 - choose v at random and a prime $e > 2$ of length higher than each m_1, \dots, m_l
 - $A := ((Z / (S^v \cdot \prod R_i^{m_i}))^{1/e})$
 - the signature is (A, e, v)
- verification: check if

$$Z = A^e \cdot S^v \cdot \prod R_i^{m_i} \quad ?$$

Issuing a certificate for values m_1, \dots, m_l

- somewhat complicated since the Issuer can learn only some attributes to be signed
- **method**: a two-party protocol to compute CL signature of the Issuer, algorithm:
 - issuer sends a challenge n_1
 - the user chooses v' at random and computes $U := S^{v'} \cdot \prod R_i^{m_i}$ apart from known attributes
 - the user creates a ZKP that U computed in this way, in particular that
 - the user knows hidden attributes
 - the user uses the same attributes as committed
 - the issuer checks the ZKP proof
 - the issuer chooses a random prime e and v''
 - the issuer computes
$$Q := Z / (U \cdot S^{v''} \cdot \prod_{\text{known } m_i} R^{m_i}) \text{ and } A := Q^{1/e}$$
 - (A, e, v'') sent to the user together with ZKP proof of correctness
 - the user computes $v := v' + v''$, checks the proof and validity of signature (A, e, v)

Presenting a credential

complicated: also involves proofs over encrypted values and the range of attributes. Some attributes may be revealed, but some stay hidden. Moreover, the certificate must not be revealed.

some details for verification of certificate without revealing it:

- values \tilde{m}_i chosen for each hidden attribute m_i , where $i \in A_{\bar{r}}$
- the user chooses r_A at random and randomizes (A, e, v) :
 - $A' := A \cdot S^{r_A}$, $v' := v - e \cdot r_A$
- so called t -values computed:
 - chosen at random: \tilde{e}, \tilde{v}'
 - $\tilde{Z} := (A')^{\tilde{e}} \cdot S^{\tilde{v}'}$
- these t -values \tilde{Z} and t values from other proofs plus some other data are hashed to get challenge c
- signatures components (s -values) are derived:
 - $\hat{e} := \tilde{e} + c \cdot e$
 - $\hat{v}' := \tilde{v}' + c \cdot v'$
 - $\hat{m}_i := \tilde{m}_i + c \cdot m_i$

Credential verification - based on recomputation of t -values and recomputing c .

\tilde{Z} recomputed as:

$$(A')^{\hat{e}} \cdot \prod_{i \in A_{\bar{r}}} R_i^{\hat{m}_i} \cdot S^{\hat{v}'} / \left(\frac{Z}{\prod_{i \in A_{\bar{r}}} R^{m_i}} \right)^c$$

(we remove some parts from Z , what is left must be raised to power c to eliminate the effect of components $c \cdot e$, $c \cdot v'$, $c \cdot m_i$ when using the exponents \hat{e} , \hat{v}' , \hat{m}_i

ranges have to be checked, etc

...

IDENTIFICATION

running wireless communication protocol may enable tracing a user.

Threats:

- explicit exchange of identifiers: an eavesdropper learns who is communicating with whom
- strong cryptographic proofs created during identification: can be misused for proving presence to the third parties

elimination of explicit identifiers:

- at each communication round Alice and Bob create random nonce (nonces) for the next round
- even more secure: if n is such a nonce, then Alice uses n' where n' is the same as n except for one bit at a random position

deniability:

- the idea is that a transcript of a communication (including the answer from the Prover created with his private key) can be simulated

consequence: a third party has no grounds to believe the communication transcript presented to him

- **wrong example:** challenge-response algorithm with digital signature:

- 1 the Verifier selects x at random and sends to the Prover
- 2 the Prover returns his signature s over x

unfortunately s can serve as a proof of the claim of the Verifier: “I have talked to Prover” if x is a signature of the Verifier or something that only could be created by the Verifier

- **good example:** static Diffie-Hellman protocol

- **good example:** Stinson-Wu for Prover with the key pair $(a, A = g^a)$

- 1 Verifier chooses x at random, computes $X := g^x$ and $Y := \text{Hash}(A^x)$ and sends X, Y to Prover
- 2 Prover computes $Z := X^a$ and aborts if $Y \neq \text{Hash}(Z)$
- 3 Prover sends Z
- 4 Verifier accepts iff $Z = A^x$

Stinson-Wu protocol

- Stinson-Wu does not create an oracle for DH Problem, Verifier must send a challenge for which somebody knows x
- it is not true that Verifier must know x :
 - Eve creates correct X, Y as well as x encrypted with $\text{Hash}(Z)$
 - Eve sends these data to Verifier in advance
 - somewhat later: Verifier sends X, Y to Prover
 - Prover computes $Z := X^a$ and aborts if $Y \neq \text{Hash}(Z)$
 - Prover sends Z
 - Verifier computes $\text{Hash}(Z)$ and uses it as a key to decrypt and derive x
 - Verifier accepts iff $Z = A^x$
 - **Verifier returns x to Eve as a proof of interaction with Prover**